

Machine Learning 1

Lecture 12.5 - Kernel Methods
Gaussian Processes - Bayesian Regression

Erik Bekkers

(Bishop 6.4.2, 6.4.3)



Regression with GP's

- ▶ We have observed $\{(\mathbf{x}_i, f_i)\}_{i=1}^N$ where we assume

$$f_i = f(\mathbf{x}_i) = y(\mathbf{x}_i) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \beta^{-1})$$

- ▶ Assume we have a GP for $y(x)$, so for any

$$\mathbf{y} = \begin{bmatrix} y(\mathbf{x}_1) \\ \vdots \\ y(\mathbf{x}_N) \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \dots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix} \right)$$

$K(\mathbf{X}, \mathbf{X})$

- ▶ Then $f(\cdot)$ is also a *GP* since $\mathbf{f} = \mathbf{y} + \boldsymbol{\varepsilon}$, and the sum of two independent random variables is also Gaussian distributed

non-parametric

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, K(\mathbf{X}, \mathbf{X}) + \beta^{-1} \mathbf{I})$$

vs

parametric

$$\underline{f} \sim \mathcal{N}(\Phi \underline{w}, \beta^{-1} \mathbf{I})$$

$\underline{w} \sim p(\underline{w})$

equivalent kernel

Predictions with GP's

- ▶ The joint distribution of test points \mathbf{f}^* (at \mathbf{X}^*) and \mathbf{f} (train points), according to our *GP*, is given by

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K}(\mathbf{X}, \mathbf{X}) + \beta^{-1} \mathbf{I} & \mathbf{K}(\mathbf{X}, \mathbf{X}^*) \\ \mathbf{K}(\mathbf{X}^*, \mathbf{X}) & \mathbf{K}(\mathbf{X}^*, \mathbf{X}^*) + \beta^{-1} \mathbf{I} \end{bmatrix} \right)$$

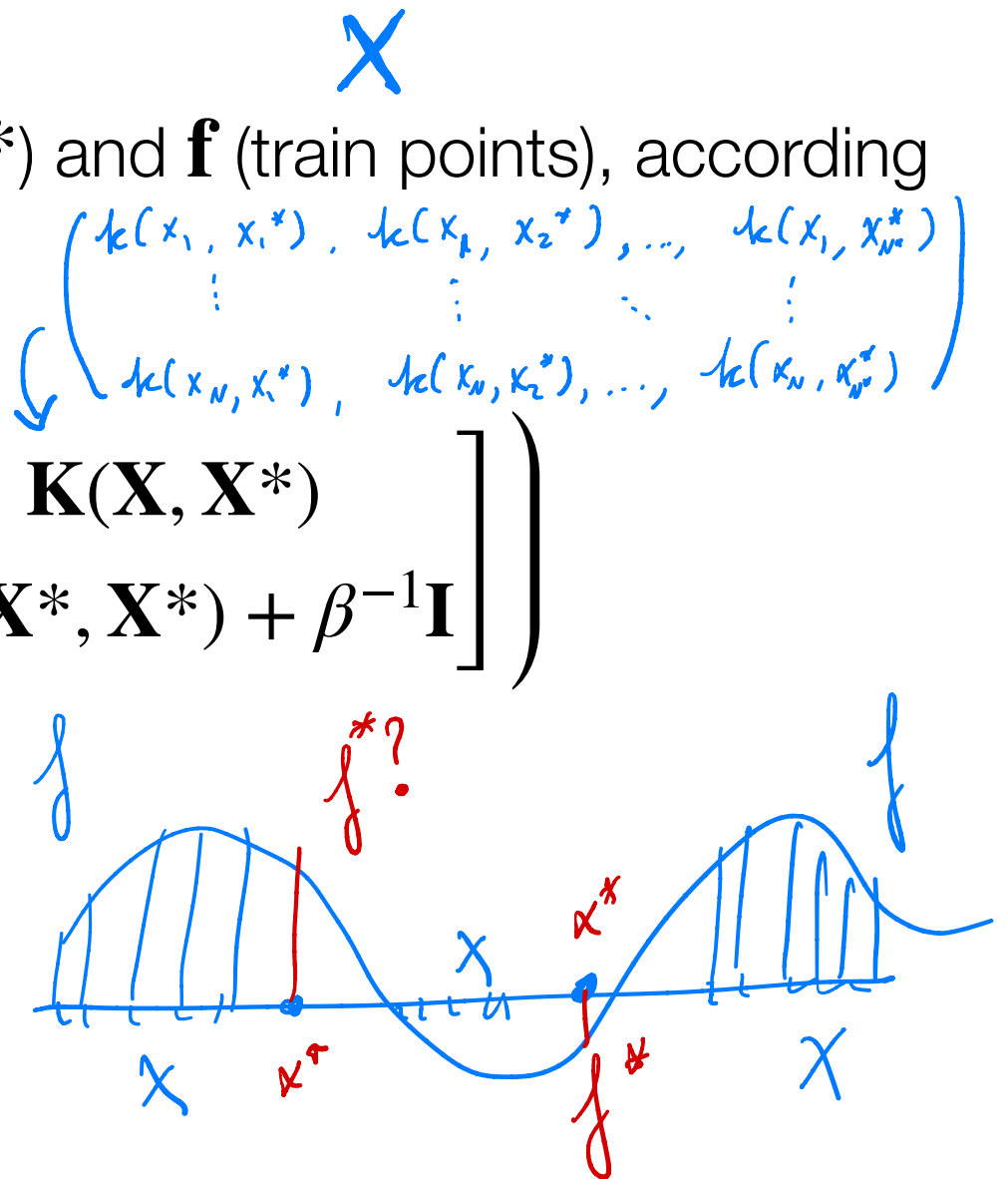
- ▶ Then *Gaussian conditioning property!*

$$p(\mathbf{f}^* | \mathbf{X}^*, \mathbf{X}, \mathbf{f}) = \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$$

with

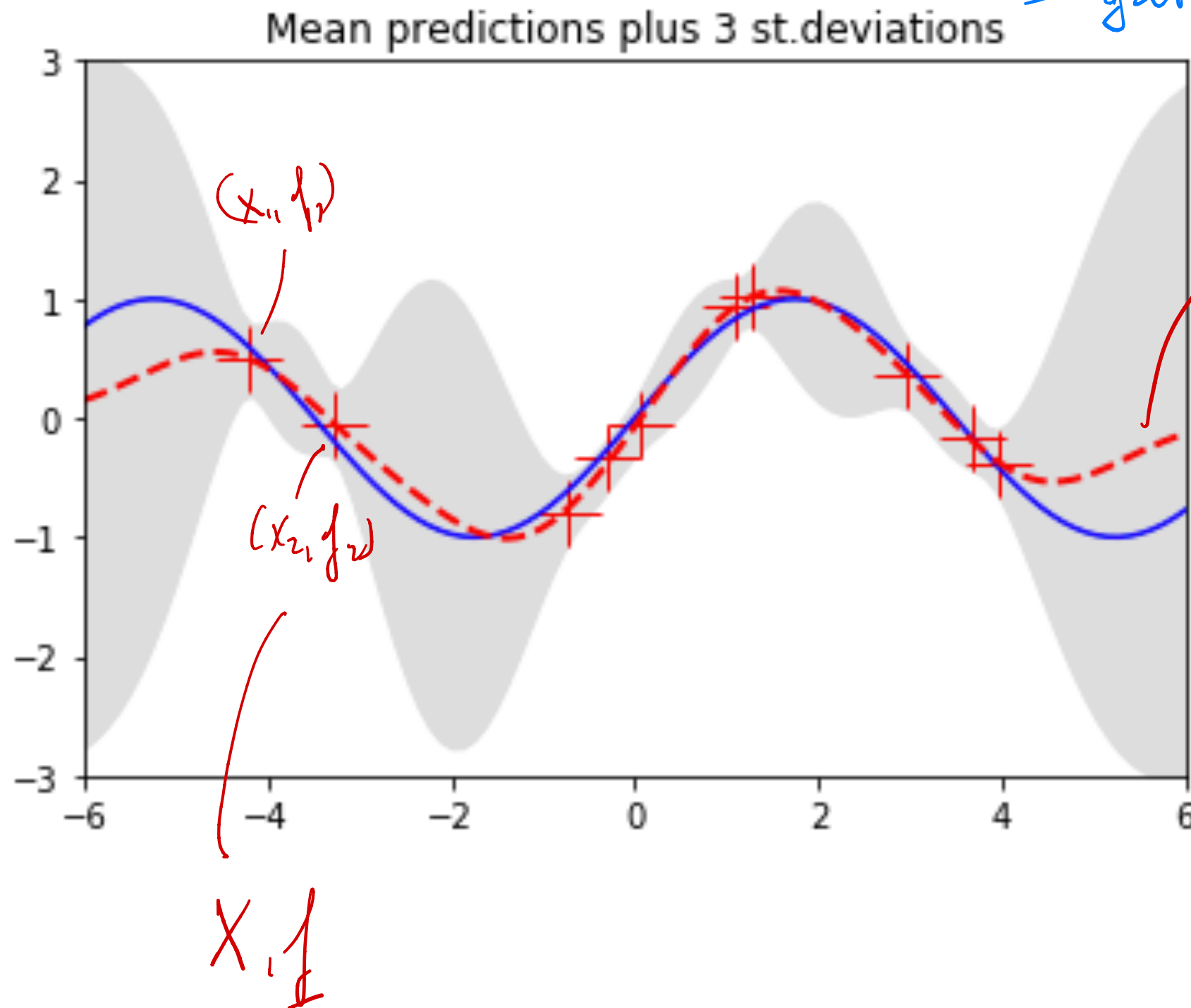
$$\boldsymbol{\mu}^* = \mathbf{K}(\mathbf{X}^*, \mathbf{X}) (\mathbf{K}(\mathbf{X}, \mathbf{X}) + \beta^{-1} \mathbf{I})^{-1} \mathbf{f}$$

$$\boldsymbol{\Sigma}^* = \mathbf{K}(\mathbf{X}^*, \mathbf{X}^*) + \beta^{-1} \mathbf{I} - \mathbf{K}(\mathbf{X}^*, \mathbf{X}) (\mathbf{K}(\mathbf{X}, \mathbf{X}) + \beta^{-1} \mathbf{I})^{-1} \mathbf{K}(\mathbf{X}, \mathbf{X}^*)$$

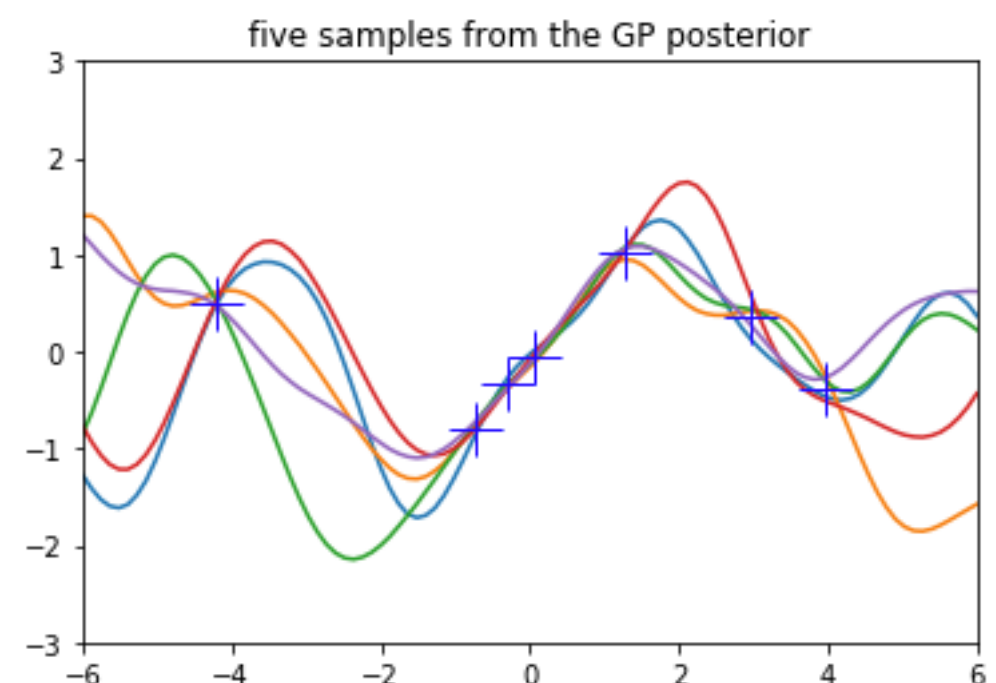
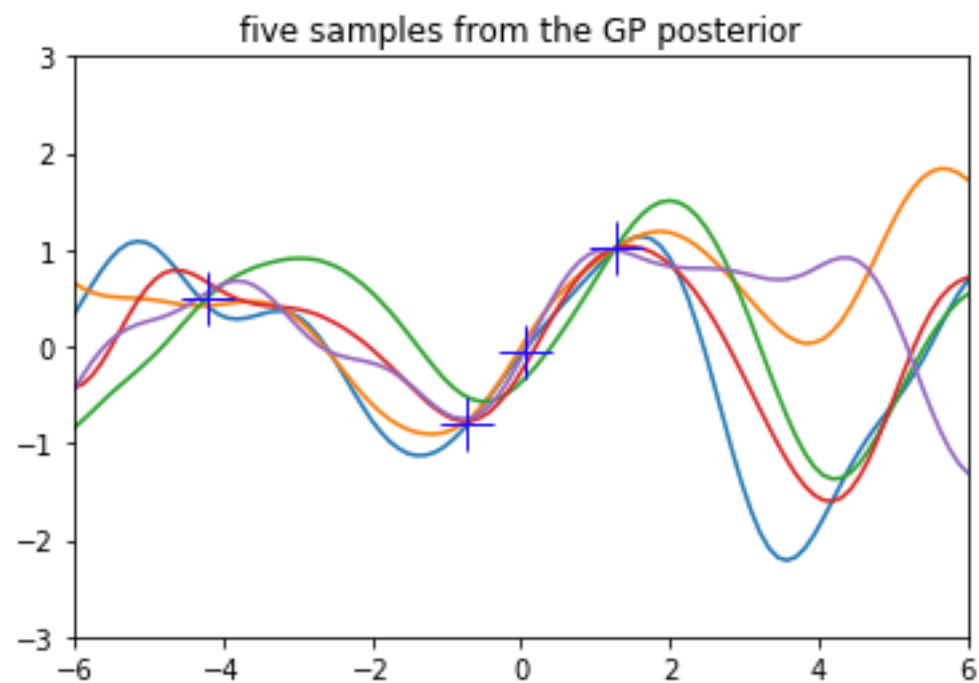
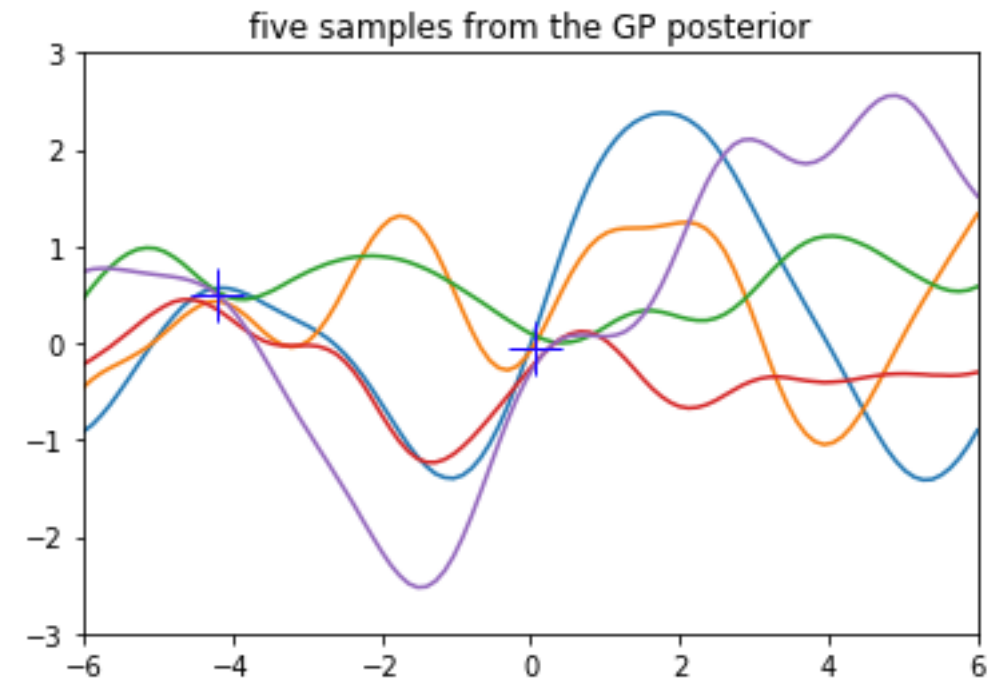
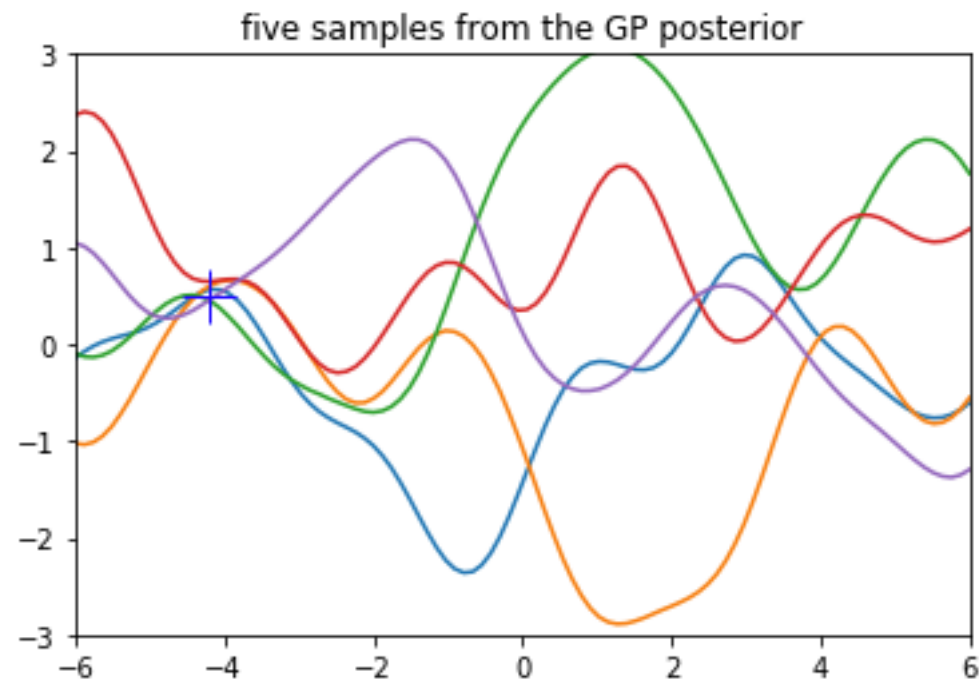


Predictions with GP's

active learning!
- identify uncertain regions
- gather more data



Drawing functions from GP posterior



How to choose kernel parameters?

- ▶ The kernel parameters $\theta_0, \theta_1, \theta_2, \theta_3$ are hyperparameters
- ▶ Simplest approach: take training observations, for which we know

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_{\theta}(\mathbf{X}, \mathbf{X})) = \frac{1}{(2\pi)^{N/2} |\mathbf{C}_{\theta}|^{1/2}} \exp \left(-\frac{1}{2} \mathbf{f}^T \mathbf{C}_{\theta}^{-1} \mathbf{f} \right)$$

$$\text{with } \mathbf{C}_{\theta}(\mathbf{X}, \mathbf{X}) = \mathbf{K}_{\theta}(\mathbf{X}, \mathbf{X}) + \beta^{-1} \mathbf{I}$$

- ▶ Make a maximum likelihood estimate

$$\max_{\theta} \ln p(\mathbf{f} | \mathbf{X}, \theta) = \max_{\theta} -\frac{1}{2} \ln |\mathbf{C}_{\theta}| - \frac{1}{2} \mathbf{f}^T \mathbf{C}_{\theta}^{-1} \mathbf{f} - \frac{N}{2} \ln 2\pi$$

- ▶ Solve numerically for θ